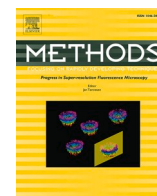




Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.



# Identification of potential pan-coronavirus therapies using a computational drug repurposing platform

Woochang Hwang<sup>a</sup>, Namshik Han<sup>a,b,\*</sup>

<sup>a</sup> Milner Therapeutics Institute, University of Cambridge, Cambridge, UK

<sup>b</sup> Cambridge Centre for AI in Medicine, Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Cambridge, UK

## ARTICLE INFO

### Keywords:

COVID-19  
SARS  
Drug repurposing  
Artificial neural network  
Drug mechanism  
Virus replication  
Coronavirus

## ABSTRACT

In the past 20 years, there have been several infectious disease outbreaks in humans for which the causative agent has been a zoonotic coronavirus. Novel infectious disease outbreaks, as illustrated by the current coronavirus disease 2019 (COVID-19) pandemic, demand a rapid response in terms of identifying effective treatments for seriously ill patients. The repurposing of approved drugs from other therapeutic areas is one of the most practical routes through which to approach this. Here, we present a systematic network-based drug repurposing methodology, which interrogates virus–human, human protein–protein and drug–protein interactome data. We identified 196 approved drugs that are appropriate for repurposing against COVID-19 and 102 approved drugs against a related coronavirus, severe acute respiratory syndrome (SARS-CoV). We constructed a protein–protein interaction (PPI) network based on disease signatures from COVID-19 and SARS multi-omics datasets. Analysis of this PPI network uncovered key pathways. Of the 196 drugs predicted to target COVID-19 related pathways, 44 (hypergeometric  $p$ -value:  $1.98e-04$ ) are already in COVID-19 clinical trials, demonstrating the validity of our approach. Using an artificial neural network, we provide information on the mechanism of action and therapeutic value for each of the identified drugs, to facilitate their rapid repurposing into clinical trials.

## 1. Introduction

Coronaviruses are pathogenic RNA viruses that cause respiratory tract infections in humans and animals. To date, seven coronaviruses have been found to cause infectious diseases in humans. Among them, severe acute respiratory syndrome coronavirus (SARS-CoV) in 2002, and COVID-19 (SARS-CoV-2) in 2019, can produce severe symptoms in infected individuals and to date have killed an estimated 750 people [1] and 4.0 million people [2], respectively. The global pandemic caused by SARS-CoV-2 is ongoing, with the number of infections and deaths steadily increasing. The initial cases of both SARS and COVID-19 are thought to have been caused by the crossing over of coronaviruses between animals and humans (zoonotic transmission) and some have predicted that as human populations continue to grow and expand into new habitats, zoonotic transmission resulting in the emergence of novel infectious diseases could become common [3].

As a result of the SARS-CoV-2 pandemic and the need for treatments to alleviate severe disease and reduce deaths, studies have been carried out at an unprecedented pace to understand this virus. As a result, virus–host interactome data that identifies directly interacting proteins

(DIPs) [4] and proteomics data that maps differentially expressed proteins (DEPs) [4] after infection have been generated and shared. Such datasets, along with and other omics data from virus studies, can form the basis of drug repurposing studies, one of the more rapid approaches for identifying potential new treatments [5,6].

Previously, we have applied a drug repurposing strategy to identify potential new treatments for COVID-19 [7]. In this study, we use a similar computational approach, but include datasets from cells infected with either SARS-CoV-2 or SARS-CoV. Our integrative drug repurposing methodology combines computational biology and machine learning approaches to analyse datasets from coronavirus studies. Our methodology uses virus–host interactions and proteins that are differentially expressed 24 h after virus infection, identified from previous experimental studies. We construct a virus-induced network to identify potential antiviral therapies and analyse mechanisms of actions of these potential therapies (Fig. 1). The application of this methodology to both COVID-19 and SARS indicates that these approaches could be used to identify potential therapies for future novel viral infectious disease outbreaks.

\* Corresponding author.

E-mail address: [nh417@cam.ac.uk](mailto:nh417@cam.ac.uk) (N. Han).

<https://doi.org/10.1016/j.ymeth.2021.11.002>

Received 1 September 2021; Received in revised form 29 October 2021; Accepted 3 November 2021

Available online 9 November 2021

1046-2023/© 2021 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 2. Materials and methods

### 2.1. Data collection

DIP and DEP data for SARS-CoV and SARS-CoV-2 were obtained from Stukalov et al. [4]. In brief, this study [4] used A549 human lung carcinoma cells infected with SARS-CoV or SARS-CoV-2 and 6, 12 and 24 h after infection mass spectrometry analyses were used to investigate virus–host interactions and changes in protein expression. We used the high-confidence virus–host interactions identified after infection as DIP for SARS-CoV (521 proteins) and SARS-CoV-2 (781 proteins). As alterations in protein expression levels were most evident at 24 h after infection (hpi; Bayesian linear model-based unadjusted two-sided  $P$ -value  $\leq 10^{-3}$ ,  $|\log_2$  fold change|  $\geq 0.5$ ), we used the published data from this time point as DEPs for SARS-CoV (188 proteins) and SARS-CoV-2 (197 proteins).

### 2.2. CIP network construction

The coronavirus induced protein (CIP) network was constructed using all shortest paths between DIP and DEP in a human protein–protein interaction network generated using the STRING database (v11.0) [8]. The main purpose of constructing the CIP network in our study was to identify COVID-19 disease associated proteins. The STRING database was selected as the PPI database because of previous evidence that it contains more comprehensive information on diverse collections of

disease-associated protein sets compared with other databases [9]. Only interactions with a confidence score of more than medium (0.7) were used. The 0.7 cut-off is high confidence level for protein–protein interaction searches in the STRING database [10,11], but also enables all potential interactions to be considered, thereby enabling as many key proteins as possible to be identified and analysed.

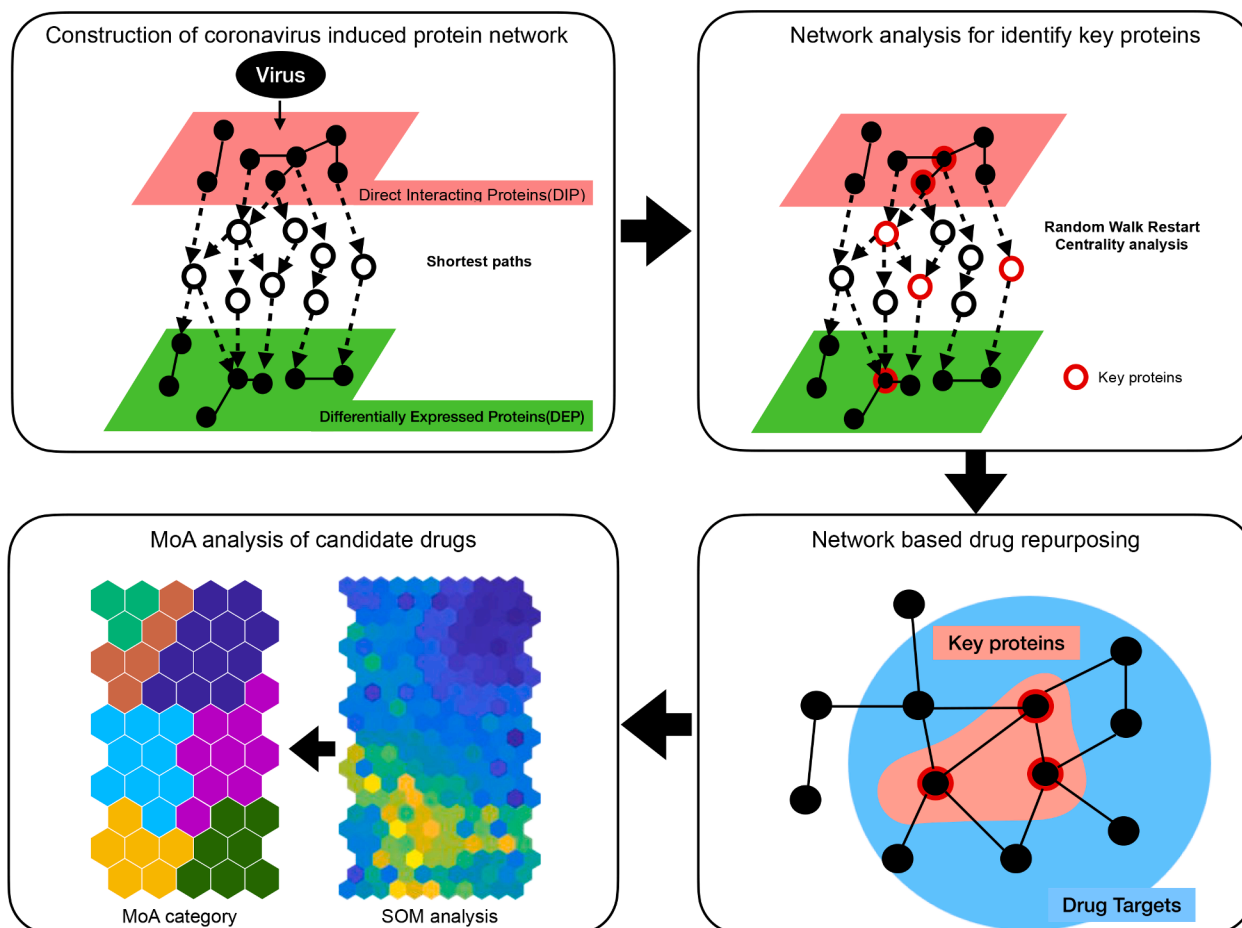
All shortest paths between all pair proteins of DIP and DEP on the human PPI network were found using Dijkstra algorithm. For the shortest path finding, we used the python package NetworkX (v2.2) [12]. All paths between DIPs and DEPs compose a “hidden layer” in the CIP network.

### 2.3. Network analysis on CIP network

Eigenvector centrality, degree centrality, betweenness centrality and random walk with restart (RWR) were utilized to identify key proteins in CIP networks. The CIP network was represented by an adjacency matrix  $A$ , where  $A_{ij} = 1$  there is an edge between nodes  $i$  and  $j$  or  $A_{ij} = 0$  otherwise. The eigenvector centrality  $x_i$  was defined as

$$\lambda x = xA \quad (1)$$

where  $x$  is an eigenvector of the adjacency matrix  $A$  with eigenvalue  $\lambda$ . If  $\lambda$  is the largest eigenvalue of the adjacency matrix  $A$ , there is a unique solution  $x$ , all centrality values are positive [13]. Degree centrality of node  $i$  was defined as



**Fig. 1.** Workflow of computational drug repurposing platform using network. The coronavirus induced protein (CIP) network was constructed by connecting through the shortest paths the directly interacting proteins (DIPs), identified using a virus–host interactome dataset, with the differentially expressed proteins (DEPs), identified using proteins that are differentially expressed 24 h after infection [2]. Network analyses are used to identify key proteins involved in coronavirus infection and then potential drug candidates are identified using a network proximity model between key proteins and drug targets. An artificial neural network is used to map the mechanism of action of the potential drugs.

$$C_D(i) = \sum_{j=1}^N A_{ij} \quad (2)$$

where  $N$  is the number of nodes in the CIP network. Betweenness centrality of a node  $i$  was defined as

$$C_B(i) = \sum_{s,t \in V} \frac{\sigma(s,t|i)}{\sigma(s,t)} \quad (3)$$

where  $V$  is the set of nodes,  $\sigma(s,t)$  is the total number of shortest paths between  $s$  and  $t$ , and  $\sigma(s,t|i)$  is the number of shortest paths between  $s$  and  $t$  passing through node  $i$ . If  $s = t$ ,  $\sigma(s,t) = 1$ , and if  $i \in s, t$ ,  $\sigma(s,t|i) = 0$ .

Fig. 2 shows how each of the centrality algorithms impact the network. Eigenvector centrality was used to identify the most influential proteins in the network (Fig. 2B). If a protein interacts with many other proteins as a hub protein, and those proteins also interact with many other proteins, the protein will have high eigenvector centrality. Degree centrality was used to identify the hub proteins in the network (Fig. 2C). Betweenness centrality measure is based on communication flow and measures how often a node occurs on all shortest paths between two nodes (Fig. 2D). Betweenness centrality was used to identify the bottleneck proteins in the network. Bottlenecks, that have many shortest paths passing through them, are connector proteins for inter-pathways with a functional property [14]. A RWR algorithm was used to see which human proteins are most affected upon SARS-CoV-2 infection. To do this, we used all DIP as the starting points of RWR. The RWR parameters were (1) a restart probability of 0.15, (2) a maximum iteration number of 100, and (3) an error tolerance of  $1e-06$ . We assigned edge betweenness centrality as an edge score on the CIP network. The RWR

calculated a score per protein in the CIP network which indicates how much a given protein was influenced by SARS-CoV-2 via DIP. The algorithms were implemented in the python package NetworkX (v 2.2) [12].

Permutation tests were performed 1000 times to identify significant proteins for each of the network centrality algorithms. In 1000 permutation tests, each test generated a random network with a preserved degree distribution of the original network, the CIP network. To generate a random network, we reconnected the edge in the CIP network and swiped the node. So, the random network in each permutation test has at least 66% of the rewired edges. Then, in the permutation test, we applied the network algorithm and obtained the cumulative results of the network algorithm. These cumulative results were used to calculate the empirical p-value of the network algorithm. We combined the four permutation test results to determine the final set of key proteins that have an empirical p-value  $\leq 0.01$  in either result (Fig. S1,2).

#### 2.4. Over-representation analysis of key proteins

To characterise key proteins of CIP networks, we performed over-representation analysis (ORA) on GO biological process associated with key proteins, gprofiler2 [15] (hypergeometric test, false discovery rate (FDR)-BH  $< 0.01$ ). We also calculate the GO-ORA similarity between SARS and COVID-19 using overlap similarity:

$$O(A,B) = \frac{|A \cap B|}{\min(|A|, |B|)} \quad (4)$$

#### 2.5. Network-based drug simulation

Approved drugs were collected from ChEMBL [16] and DrugBank

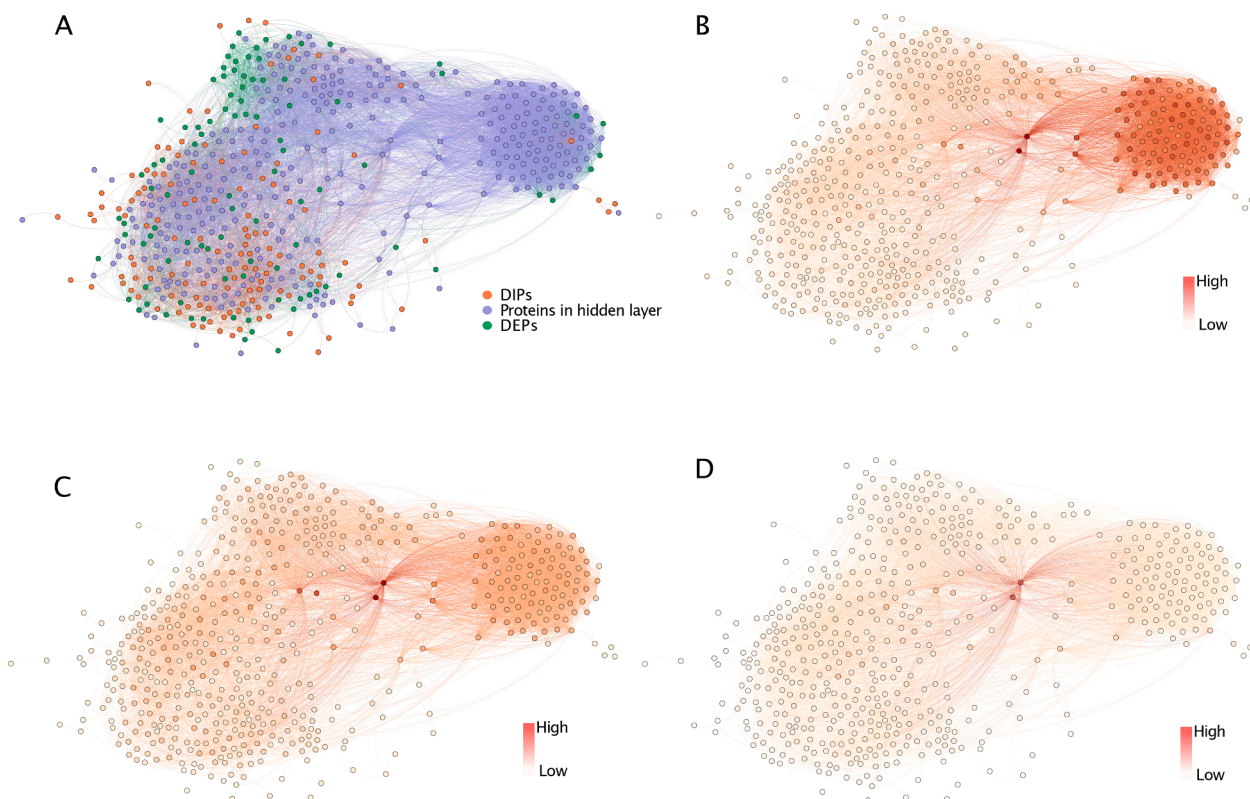


Fig. 2. Proteins in the hidden layer show significant centrality score in eigenvector, degree and betweenness centrality. Nodes in the network are key proteins of COVID-19 network. A significant node has a higher centrality score. Edges are the interactions between key proteins. Edge colours represent average centrality score of two interacting proteins. (A) Shows a subnetwork of the COVID-19 network composed of key proteins. (B) Illustrates the result of applying eigenvector centrality to this subnetwork. (C) Shows the result of applying degree centrality to this subnetwork. (D) Illustrates the result of applying betweenness centrality to this subnetwork.

[17]. Drug-target interaction information was collected from DrugBank (v 5.1) [17], STITCH (v 5.0, confidence score >0.9) [18] and Cheng, et al. [19].

*In-silico* network-based proximity analysis was conducted for key proteins from the CIP network for SARS and COVID-19. Given  $K$ , the set of key proteins from CIP networks, and  $T$ , the set of drug targets, the network proximity (Eq. (5)) of  $K$  with the target set of  $T$  of each approved drug where  $d(k, t)$  the shortest path length between nodes  $k \in K$  and  $t \in T$  in the human PPIs [19] was executed. Closest distance measure was used to calculate the distance between a given drug's targets and key proteins in the CIP network, because it previously showed best performance in drug-disease pair prediction in the study of Guney et al. [20].

$$d_c(K, T) = \frac{1}{\|T\|} \sum_{t \in T} \min_{k \in K} d(k, t) \quad (5)$$

To assess the significance of the distance between a key protein in the CIP network and a drug  $d_c(K, T)$ , the distance was converted to Z-score based on permutation tests by using

$$z_c(K, T) = \frac{d_c(K, T) - \mu_{d_c(K, T)}}{\sigma_{d_c(K, T)}} \quad (6)$$

The permutation tests were repeated 1000 times, each time with two randomly selected gene sets. There are few high degree nodes due to the scale-free network of the human protein-protein interaction network. To avoid repetitive selection of the same high degree nodes during random selection, we used a binning approach with at least 100 nodes in a bin. In the binning approach, nodes in same bin have a similar node degree to maintain node degree distribution for random selection. When we randomly select a set of genes, we perform a random selection among proteins from all bins to verify that the minimum node degree was less than the minimum node degree of the selected gene set and the maximum node degree was greater than the maximum node degree of the selected gene set. The corresponding p-value was calculated based on the permutation test results. The Z-score captures the statistical significance of the drug target-disease protein distance compared with the respective random expectation. Drug to disease associations with a Z-score of less than  $-2$  were considered significantly proximal [20]. We selected drugs with a Z-score of less than  $-2$  as potential drugs for SARS or COVID-19.

## 2.6. Drug-pathway association matrix

To understand mechanisms of action (MoA) for our identified drugs, we conducted the Reactome pathway over-representation analysis (ORA) of target proteins of drugs using R (v 3.5.2) package, gprofiler2 [15] (hypergeometric test, FDR-BH < 0.01). Reactome pathways (the version on 15/05/2020) was used for pathway ORA because it is most actively updated public database of human pathways [21]. Pathway ORA was conducted by extending not only the target protein of a drug but also adjacent proteins to identify a drug-related pathway — through utilizing the advantages of network analysis that considers not only a single protein but also associated neighbour proteins.

Significantly enriched biological pathways of drug targets for each of the identified drugs were integrated, resulting in 865 pathways (FDR-BH <  $1.0e-13$ ) for SARS-CoV and 1008 pathways (FDR-BH <  $1.0e-12$ ) for SARS-CoV-2. To include all drug-associated pathways, we integrated pathways while reducing the FDR. The Reactome pathway has a hierarchical structure among pathways, whereby the lower hierarchy pathway is more specific than the higher hierarchy pathway. The parent pathway semantically includes the child pathways. In the process of integrating the enriched pathways per drug, we filtered the topmost pathway (e.g. Metabolism, Apoptosis, Cell cycle) in pathway hierarchy of Reactome. The lowest pathways are merged to its parent pathways. For example, “Regulation of PTEN stability and activity” and

“Regulation of PTEN gene transcription” are merged to “PIP3 activates AKT signalling”. We identified 219 key pathways associated with identified drugs of SARS-CoV and 249 key pathways associated with identified drugs of SARS-CoV-2.

Based on these identifications, a matrix containing F1 scores of the identified drugs and the key pathways was generated for drug-pathway association. The Reactome pathway ORA for identified drugs using gprofiler2 provides enrichment p-values, precision and recall information that were used to produce the F1 scores.

$$F1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (7)$$

The meaning of precision here is the proportion of drug targets that are annotated to the pathway. The meaning of recall here is the proportion of the pathway gene set that the drug targets recover. The pathway to which the largest number of drug target proteins belong has the highest precision value. The pathway with the greatest intersection of pathway proteins and target proteins has the highest recall value. In other words, the pathway with the highest F1 score in the drug-pathway associations is the pathway to which the drug's target protein associates most strongly and the pathway with the largest intersection between the target proteins. For example, the number of target proteins for dexamethasone is 12. The number of 'Interleukin-4 and Interleukin-13 signaling' pathway proteins is 112. The number of intersections between the target protein of dexamethasone and the 'Interleukin-4 and Interleukin-13 signaling' pathway protein is 4. So, the precision is  $4/12 = 0.33$  and the recall value is  $4/112 = 0.0357$ . Thus, the F1 score is 0.064. The number of 'Cytokine Signaling in Immune system' pathway proteins is 849, and the number of intersections between the target protein of dexamethasone and the 'Cytokine Signaling in Immune system' pathway protein is 4. The precision is the same as 0.33 for 'Interleukin-4 and Interleukin-13 signaling', but the recall value is  $4/849 = 0.0047$ . Thus, the F1 score is 0.0092, which is lower than the 'Interleukin-4 and Interleukin-13 signaling'. As such, the F1 score complements the imbalance between the pathway protein and the target protein.

## 2.7. SOM analysis

Self-Organizing Map (SOM) [22] was used to analyse the MoA of the candidate drugs from our analysis. SOM has descriptive ability and hence advantages in visual concept detection. Thus, it was useful to directly compare the SOM component heatmaps of the 249 pathways for COVID-19 and 219 pathways for SARS. SOM also has the advantage of dimensional reduction to allow a more appropriate clustering result. SOM was used to calculate the low-dimensional abstractions, which are then clustered using k-means. This two-phase approach increases the efficiency of k-means clustering with a relatively small number of samples, a known limitation in hierarchical clustering algorithms including k-means. Another advantage of SOM is noise reduction because the SOM abstractions are less sensitive to random variations than the input data. In addition, SOM offers a systematic arrangement of the candidate drugs to each neuron (or node) and hence to pathway clusters.

The data used in training was the F1 score matrix for Drug-Pathway associations. After the SOM training, the Unified distance matrix (U-matrix) of the trained unsupervised SOM contains the vector norms between the neighbouring SOM nodes and shows data density in input space. Each sub-unit is coloured according to distance between corresponding data vectors of neighbour units. Low distances areas (dark blue) have high data density (clusters) (Fig. S3,4). Davies-Bouldin index (DBI) [23] was calculated based on the U-matrix to determine the optimal number of clusters for the k-means algorithm. The k-means algorithm was then used to cluster pathways for COVID-19 and SARS. The lowest DB index value occurred at 8 clusters for COVID-19 and 7 clusters for SARS, and thus we decided to map the 249 key pathways into 8

clusters and 219 key pathways into 7 clusters based on the SOM component maps of the key pathways (Fig. S3,4). The 8 clusters for COVID19 and the 7 clusters for SARS mapped to three MoA categories based on the biological function of pathways belonging to each cluster. The mapping result of key pathways to clusters and three MoA groups is available in [Supplementary Table 1](#). The detailed information of the labelled SOM neurons and the candidate drugs is available in [Supplementary Table 1](#). The SOM Toolbox package [24] for MatLab was used for this analysis with default settings and parameters.

### 3. Results

#### 3.1. Construction of coronavirus-induced protein interactome network

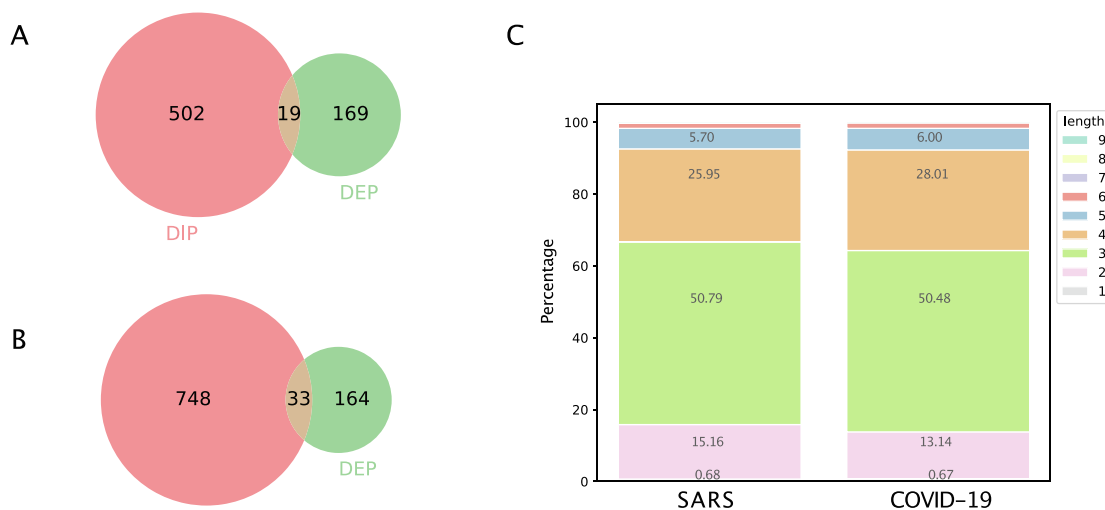
To understand mechanisms underlying SARS and COVID-19, we conducted a comprehensive analysis of the protein sets implicated in each virus infection using our workflow for data integration and network construction. To this end, we hypothesized the DIPs are the “cause” of coronavirus infection and that DEPs are the “consequence”, and there is a “hidden layer” between DIPs and DEPs where key pathways allow DIPs to control DEPs. To construct the hidden layer, we connected all pairs of DIPs and DEPs by using human protein–protein physical interaction information. We then identified all possible shortest paths between 521 DIPs and 188 DEPs in the CIP network for SARS and between 775 DIPs and 197 DEPs in the CIP network for COVID-19. As a result, we created a CIP for SARS and COVID-19. There are 10,256 proteins and 222,423 interactions in the SARS network and 11,006 proteins and 247,788 interactions in the COVID-19 network. In both SARS and COVID-19, the overlap between DIPs and DEPs is only 3% and 4%, respectively (Fig. 3A). More than 99% of pairs of DIPs and DEPs are connected through at least one protein in the hidden layer (Fig. 3B). This suggests that the hidden layers are important to understand key proteins and pathways implicated in coronavirus infection.

#### 3.2. Network analysis for identifying key proteins and pathways

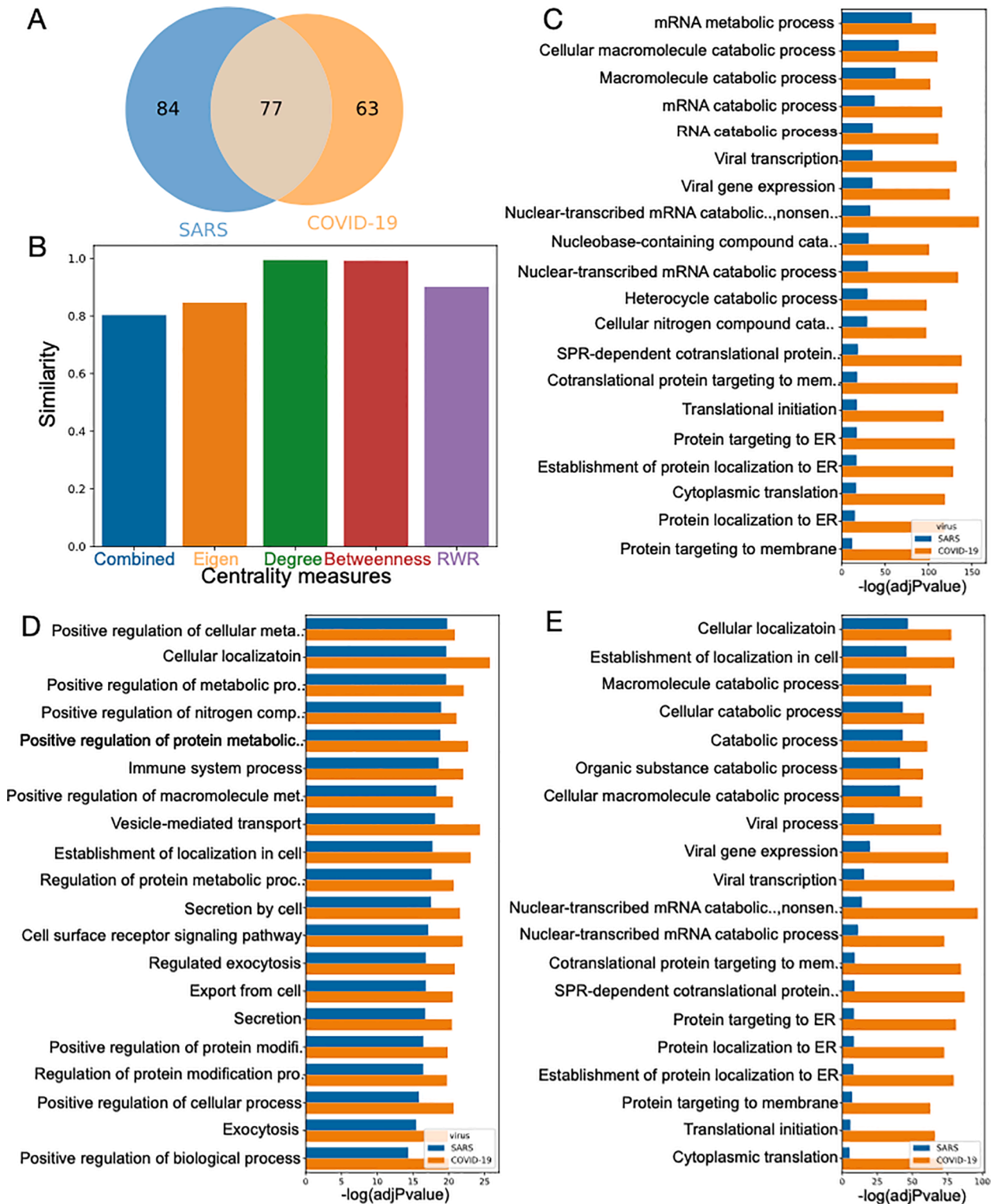
In a PPI network, proteins are key components of pathways and can significantly influence pathways as potential drug targets. We applied multiple network algorithms, including eigenvector centrality, degree centrality, betweenness centrality, and random walk with restart (RWR), to identify coronavirus related key proteins and key pathways in

the CIP network. We performed 1000 permutation tests for each network algorithm then selected proteins with empirical  $p$  values  $< 0.01$  as key proteins. We then performed ORA on Gene Ontology (GO) Biological Process (BP) to assess whether the key proteins functionally relate to the biology of coronavirus infection.

Eigenvector centrality identifies the most influential proteins in the network. It considers ‘how many proteins interact with’ and ‘what proteins interact with’ simultaneously. Although the key proteins identified by eigenvector centrality are about 55% identical between SARS and COVID-19, the enriched GO BP is 84.6% identical among them (Fig. 4A, B). The top enriched GO BPs of key proteins identified by eigenvector centrality are viral replication related terms such as mRNA metabolic processes, viral transcription, and viral gene expression [25,26] (Fig. 4C). ORA adjective  $p$ -value of GO BPs in which the key proteins of COVID-19 are overrepresented are higher than those of SARS. This suggests that the key proteins of COVID-19 are more strongly associated with viral replication-related functions than key proteins in SARS. Degree centrality identifies the hub proteins in the network. It considers only ‘how many proteins interact with’ in a distal sub-network. The ORA of the key proteins identified by degree centrality revealed mRNA metabolic processes and cell surface receptor signalling pathways which are related to viral replication, viral processes, and virus-host cell interactions (Fig. S5A) [27]. Betweenness centrality identifies the proteins that play a bridge role in a network. It considers how much a given protein is in-between others. The key proteins identified by betweenness centrality are enriched in regulatory processes, such as regulation of metabolic processes, regulation of protein metabolic processes; and communication-related biological processes, such as immune system processes. It is also associated with subcellular localization-related processes such as vesicle-mediated transport (Fig. 4D). RWR measures the importance of each protein in a network, based on the number of incoming interactions and the importance of the corresponding starting proteins, which are DIPs. RWR identifies proteins affected the most upon coronavirus infection. Key proteins identified by RWR were overrepresented in immune responses and cell-to-cell communication such as exocytosis, export from the cell and secretion (Fig. S5B). Taken together, the key proteins of CIP networks were overrepresented in virus replication-related biological processes, such as the viral process, viral transcription, and protein localization to the endoplasmic reticulum (ER). This is particularly relevant as the ER is the favourite intracellular niche for viral replication and assembly (Fig. 4E)



**Fig. 3.** Characteristics of the CIP network. (A) Ven diagram of directly interacting proteins (DIPs) and differentially expressed proteins (DEPs) for SARS. DIPs and DEPs were identified from cells infected with SARS-CoV or SARS-CoV-2 [2]. The overlap between DIPs and DEPs is only 3%. (B) Ven diagram of DIPs and DEPs for COVID-19. The overlap between DIPs and DEPs is 4%. (C) The stacked bar chart shows percentage of the length of shortest paths between DIPs and DEPs at SARS and COVID-19. There are 10,256 proteins and 222,423 interactions in the SARS network and 11,006 proteins and 247,788 interactions in the COVID-19 network.



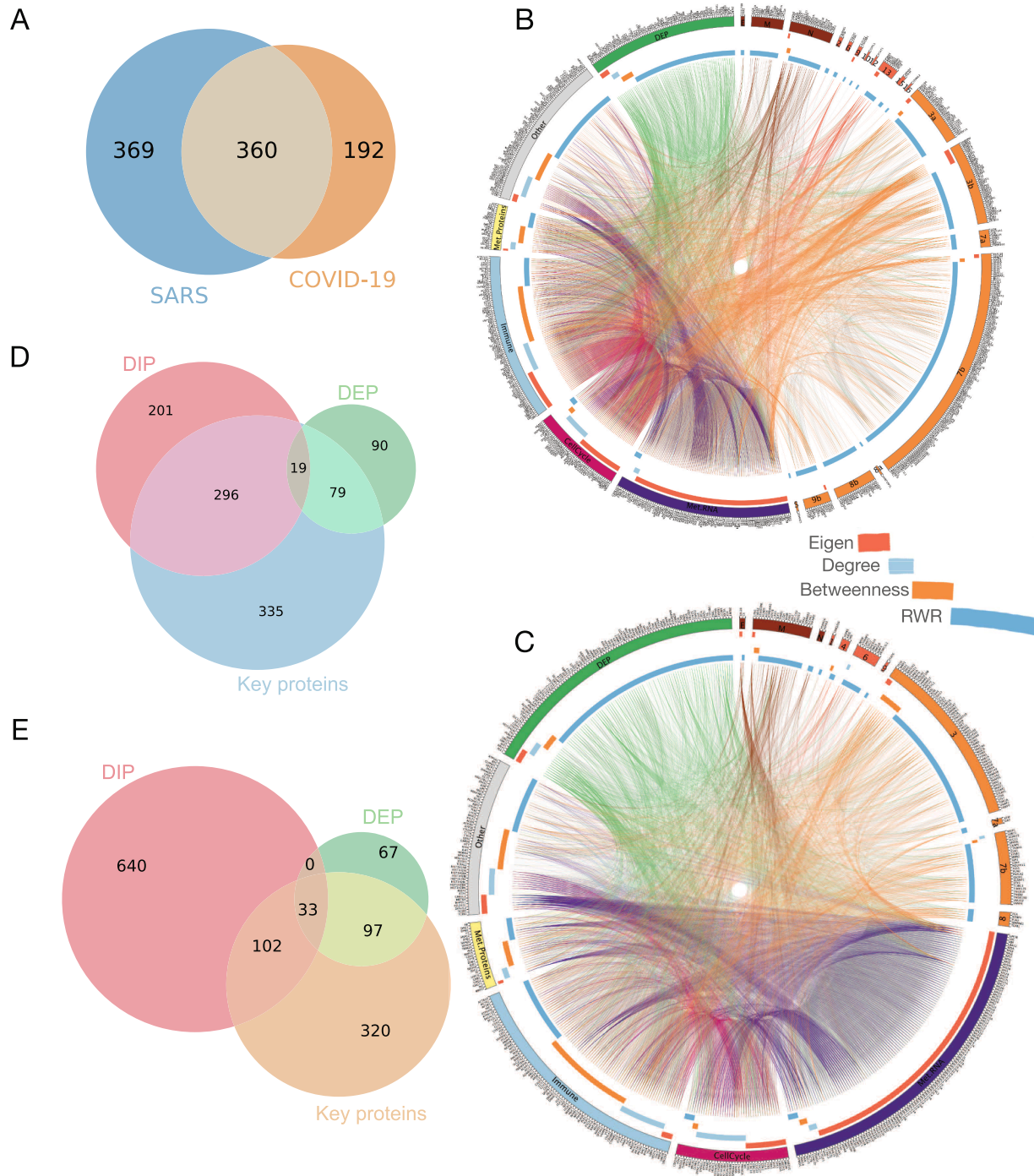
**Fig. 4.** Key proteins of SARS and COVID-19 show similar enriched biological functions. (A) Venn diagram of key proteins for both viruses, identified using an eigenvector centrality approach (B) A bar chart of the ORA similarity between SARS and COVID-19, calculated using different network centrality measures (C) A bar chart showing enriched biological functions of the key proteins by eigenvector centrality for SARS (blue) and COVID-19 (orange) (D) A bar chart showing enriched biological functions of the key proteins by betweenness centrality for SARS (blue) and COVID-19 (orange). (E) A bar chart showing enriched biological functions of combined key proteins by all network centrality measures for SARS (blue) and COVID-19 (orange). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

[28].

3.3. Comparative analysis of disease mechanisms of SARS and COVID-19

SARS and COVID-19 share 360 key proteins but also have 369 and

192 disease-specific key proteins respectively (Fig. 5A). To compare the disease mechanisms of SARS and COVID-19, we sought the most relevant biological pathways that have previously been described for SARS-CoV-2: viral replication (VR) and immune response (IR) [29]. At the highest hierarchical level in the Reactome pathway database, ‘Immune



**Fig. 5.** Comparative analysis of disease mechanisms of SARS and COVID-19. (A) A Venn diagram of key proteins in SARS and COVID-19 networks. (B) Circos plot depicting interactions between key proteins from the DIPs, DEPs and hidden layer proteins in SARS. DIPs were subdivided into the coronavirus proteins. Hidden layer proteins were subdivided into five groups: four Reactome pathways related to viral replication and immune response and one ‘Others’ that do not belong to four pathways. All key proteins are identified by eigenvector, betweenness, degree and random walk restart (RWR) network algorithms. (C) Circos plot depicting interactions between key proteins from the DIPs, DEPs and hidden layer proteins in COVID-19. (D) Venn diagram of DIPs, DEPs and key proteins in SARS. (E) Venn diagram of DIPs, DEPs and key proteins in COVID-19.



System' (SARS FDR-BH: 2.05e–23; COVID-19 FDR-BH: 4.96e–19) [25] was identified for the immune response related biological pathway. The 'Metabolism of RNA' (SARS FDR-BH: 3.11e–37; COVID-19 FDR-BH: 3.53e–58) [25,26] 'Metabolism of Protein' (SARS FDR-BH: 3.45e–15; COVID-19 FDR-BH: 4.56e–42) [30] and 'Cell Cycle' (SARS FDR-BH: 6.66e–20; COVID-19, FDR-BH: 1.59e–15) [31] were identified for viral replication related biological pathways. The key proteins belonging to these four pathways were assigned to the four band groups under the 'Hidden Layer' in Fig. 5B and C. The key proteins that did not belong to any of the four pathways were assigned to 'Other'. Among 335 key proteins (Fig. 5D) in the hidden layer of the SARS CIP network, 165 key proteins (eigen: 113, degree: 22, betweenness: 13, RWR: 17) are related to virus replication, and 88 key proteins (eigen: 22, degree: 16, betweenness: 33, RWR: 17) are related to immune response. In addition, 82 key proteins in 'Other' are overrepresented in "Fatty acid metabolism process" (FDR-BH: 4.77e–2) that are related to virus replication [32]. Among 320 key proteins (Fig. 5E) in the hidden layer of the COVID-19 CIP network, 178 key proteins (eigen: 120, degree: 26, betweenness: 15, RWR: 17) are related to virus replication, and 86 key proteins (eigen: 4, degree: 18, betweenness: 37, RWR: 27) are related to immune response. In addition, 56 key proteins in 'Other' are also overrepresented in "Fatty acid metabolism process" (FDR-BH: 1.40e–2). This might therefore potentially indicate differences in the regulation of viral replication between COVID-19 and SARS.

Comparing the four network algorithms in both diseases, we found that each algorithm identified key proteins in both diseases that have highly similar biological functions. For example, eigen centrality identified key proteins that are enriched in virus replication-related pathways. By contrast, betweenness centrality identified key proteins that are enriched in immune response related pathways. Our betweenness centrality analysis of the COVID-19 CIP network identified key viral proteins that are located in M, ORF3 and ORF7b (Fig. 5C). Interestingly, Stukalow et al. [4] reported that the M, ORF3 and ORF7b proteins of SARS-CoV-2 modulate the innate immune system. These data indicate that SARS-CoV-2 might manipulate the innate immune system to promote virus replication [4].

### 3.4. Network-based drug repurposing for SARS and COVID-19

We next sought to identify approved drugs that are able to target the identified key proteins in the CIP pathway. We conducted an *in silico* network-based proximity measure analysis [20] on the key proteins for both diseases. We collected 1917 approved drugs from publicly available databases (ChEMBL [16] and DrugBank [17]). This network-based drug repurposing analysis identified 102 drugs and 196 drugs (Supplementary Table 2) that are predicted to target the key proteins of SARS and COVID-19 CIP network, respectively. We then checked the Anatomical Therapeutic Chemical (ATC) code (84 drugs for SARS and 155 drugs for COVID-19 only were available) to determine the therapeutic areas for which the predicted drugs are normally used. The top clinical applications are cancer, sex hormone modulation, lipid signaling modulation, inflammatory/rheumatic disease, bacterial disease and immune modulation (Fig. S6).

Among the 196 identified drugs for COVID-19, 44 (22%) are currently in COVID-19 clinical trials [33] (Supplementary Table 2). To determine the significance of this finding, we asked what the likelihood would be of this number of drugs being identified as hits by chance. Given that only 13% of approved drugs (253 of 1917) are in COVID-19 clinical trials [33], the likelihood that our approach identified 44 drugs by chance is low. A hypergeometric test for the probability of 22% of our 196 drugs being in clinical trials returned a p-value of 1.98e–04, demonstrating the reliability and validity of our computational approaches. The full list of 196 approved drugs along with their detailed information is shown in Supplementary Table 2.

### 3.5. Artificial neural network for categorizing mechanisms of actions

We next wanted to discover the MoA underlying the identified drugs. In particular, we wanted to cluster the pathways and mechanisms in order to better evaluate their potential effect and utility. For SARS, an initial pathway ORA performed on proteins targeted by 102 drugs identified a set of 219 key pathways. For COVID-19, 249 key pathways have been identified as drug associated pathways (see Methods). We then calculated the precision and recall of the over representation analysis to produce an F1 score that is the measure of the ORA accuracy (see Methods). By calculating F1 scores per drug-pathway association, we generated a F1 score matrix. To investigate the MoA for identified drugs, we used SOM, a type of artificial neural network, to analyse the relationship between drugs and key pathways (termed drug-pathway association) for SARS and COVID-19. SOM, an unsupervised learning model, clusters drugs and pathways with similar patterns on the F1 matrix without prior knowledge of drugs and pathways. Then, the pathway cluster was mapped to the two therapeutically important MoAs related to coronavirus infection (viral replication and immune response). We used the pathway clusters to map the identified drugs to these two categories of MoA. The detailed information of the pathway clusters and MoAs per drug is provided in Supplementary Table 2.

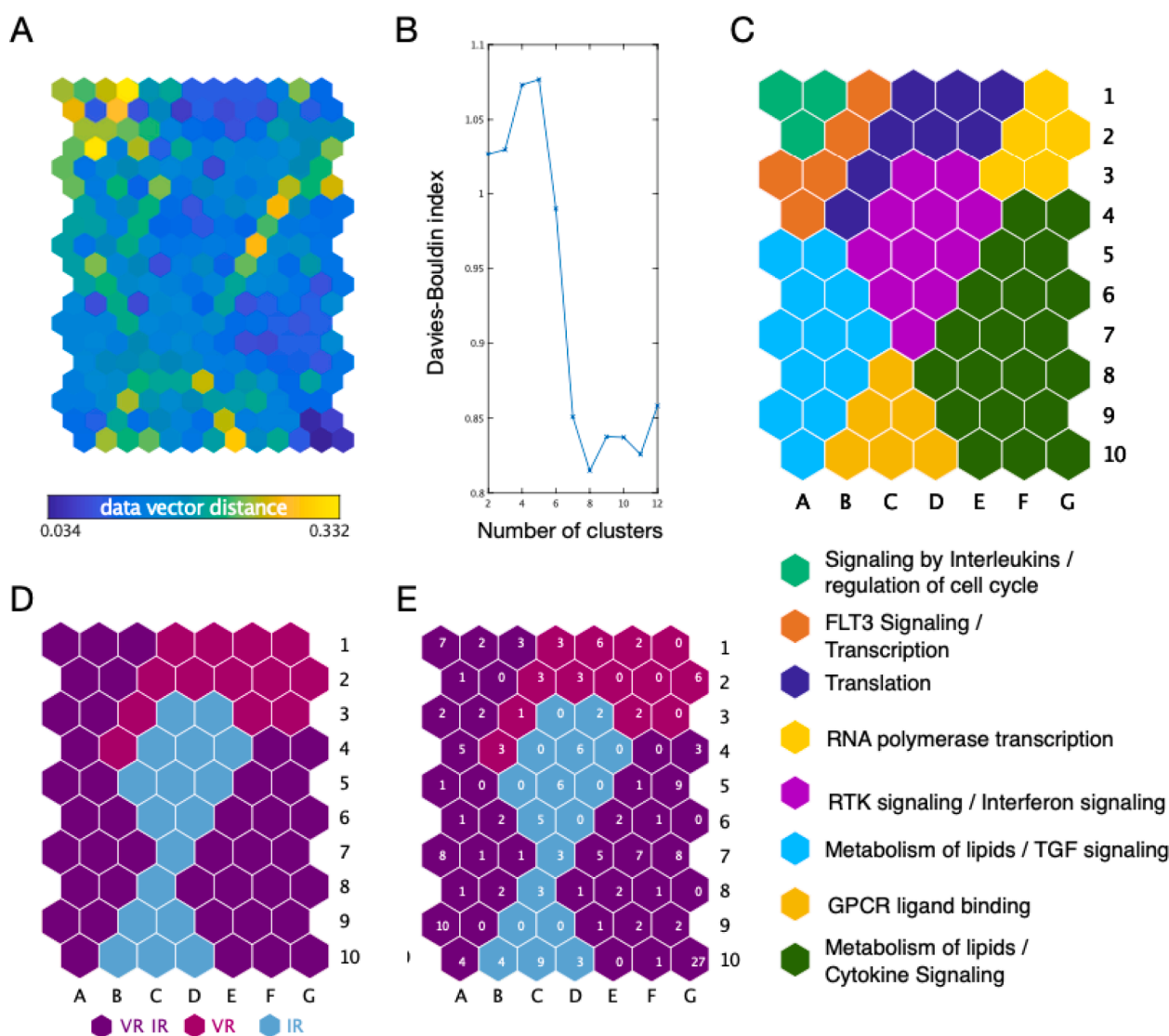
#### 3.5.1. MOAs of identified drugs for COVID-19

For COVID-19, we generated a F1 score matrix (for 196 drugs and 249 pathways that drug target proteins are overrepresented) then we investigated the MoA of identified drugs using SOM.

After training SOM with a F1 score matrix (Supplementary Table 3), we used SOM to generate 249 component plane heatmaps in order to characterise each of the 249 key pathways. Each heatmap represents the intensity patterns of a pathway. To summarise the correlation of 249 heatmaps, a U-matrix between the neighbour neurons (hexagons) was also calculated [34] (Fig. 6A). Each hexagon, which clustered drugs are located in, is a unique neuron or 'node' of the SOM. To allow direct comparison between heatmaps (pathways), the hexagons (neurons) have the same position across all heatmaps.

Next, the 249 key pathways were separated into eight clusters by a k-mean clustering algorithm with Davies-Bouldin (DB) index, which provides the optimal number of clusters (Fig. 5B). The eight clusters were "Signaling by Interleukins or regulation of cell cycle", "FLT3 signaling or transcription", "Translation", "RNA polymerase transcription", "RTK signaling or Interferon signaling", "Metabolism of lipids or TGF signaling", "GPCR ligand binding", "Metabolism of lipids or Cytokine Signaling" (Fig. 6C). The eight pathway clusters were mapped to therapeutic MoA related to COVID-19 through pathway analysis. We first sought those that are potentially major therapeutic strategies to overcome SARS-CoV-2. Novel ORFs encoded in SARS-CoV-2 are associated with viral replication or host immune response modulation [35] and treatment strategies for COVID-19 are largely divided into blocking the replication of SARS-CoV-2 and modulating the immune response [36]. Therefore, we mapped the eight pathway clusters into two categories: viral replication and immune response, using literature-driven biological supporting evidence based on the Reactome pathway hierarchy information (Table 1, Fig. 6D). For example, Cluster 3, which contains 52 pathways, is related to RNA polymerase transcription which plays an essential role in RNA viral replication [37]. Cluster 4 with 48 pathways is involved in RTK signalling or interferon signalling that are related to the host immune response during virus infection [38].

Finally, the SOM assigned 196 drugs into each hexagon to provide an indication of MoA per drug (the number of drugs per hexagon is shown in Fig. 6E). Notably, nine out of the 44 drugs that are currently in COVID-19 clinical trials [33] mapped to the 'viral replication' MoA category, 30 drugs mapped to 'viral replication or immune response', and five drugs mapped to 'immune response' (Fig. 6E, Supplementary Table 2). We then considered the mechanistic roles and connections for the 196 drugs and their target proteins that had been mapped into the



**Fig. 6.** Predicted mechanism of actions of 196 candidate drugs for COVID-19. We used SOM, a type of artificial neural network, to analyse the relationship between drugs and key pathways (termed drug–pathway association). SOM is an unsupervised learning model, which clusters drugs and pathways with similar patterns without prior knowledge of drugs and pathways. (A) U-matrix is shown of the trained unsupervised SOM used to analyze the relationship between the 196 drugs and the 249 key pathways. (B) 249 key pathways were separated into eight clusters by a k-mean clustering algorithm with Davies-Bouldin (DB) index. (C) SOM component maps (Fig. S3) are used to map the 249 pathways to the eight clusters, which are listed below the map with their associated colour. (D) The eight clusters are mapped onto three main therapeutic mechanisms of action using pathway analysis based on the Reactome pathway hierarchy information. (E) The SOM then assigns the 196 drugs to each relevant hexagon (neuron) to indicate a relevant MoA per drug.

**Table 1**  
MoA of identified drugs for COVID-19.

Cluster	Pathway	MoA category	Drugs	Drugs in clinical trials
C1	Signaling by Interleukins/ regulatin of cell cycle	VR, IR	13	4
C2	FLT3 Signaling/Transcription	VR, IR	13	10
C3	Translation	VR	17	10
C4	RNA polymerase transcription	VR	8	0
C5	RTK signaling, Interferon signaling	IR	22	2
C6	Metabolism of lipids/TGF signaling	VR, IR	35	9
C7	GPCR ligand binding/Platelet Aggregation	IR	16	3
C8	Metabolism of lipids/ Cytokine Signaling	VR, IR	72	6

eight pathway clusters. For example, Sirolimus and 16 other drugs belonging to Cluster 3 are associated with the translation function that is related to virus replication, and thus the MoA category as a COVID-19 therapy is virus replication (Supplementary Table 2). 10 out of the 16 drugs in Cluster 3 are in a COVID-19 clinical trial. Predicted drugs in Cluster 3 that are not used for cancer therapy and are not in clinical trials were then assessed for their association with viral replication through literature searches. We found literature documenting that Zinc deficiency has been confirmed in COVID-19 patients, and SARS-CoV-2 replication was reduced when zinc was administered [39]. Ademetionine has also been identified as a repurposing candidate in several studies [7,40]. And Diethylstilbestrol and adenosine have been shown to be effective in regulating the replication of other viruses [41,42]. Thus, these previous studies lend further support for the potential relevance of these drugs in Cluster.

### 3.5.2. MoA of identified drugs for SARS

We applied the same approach as outlined for COVID-19 to identify

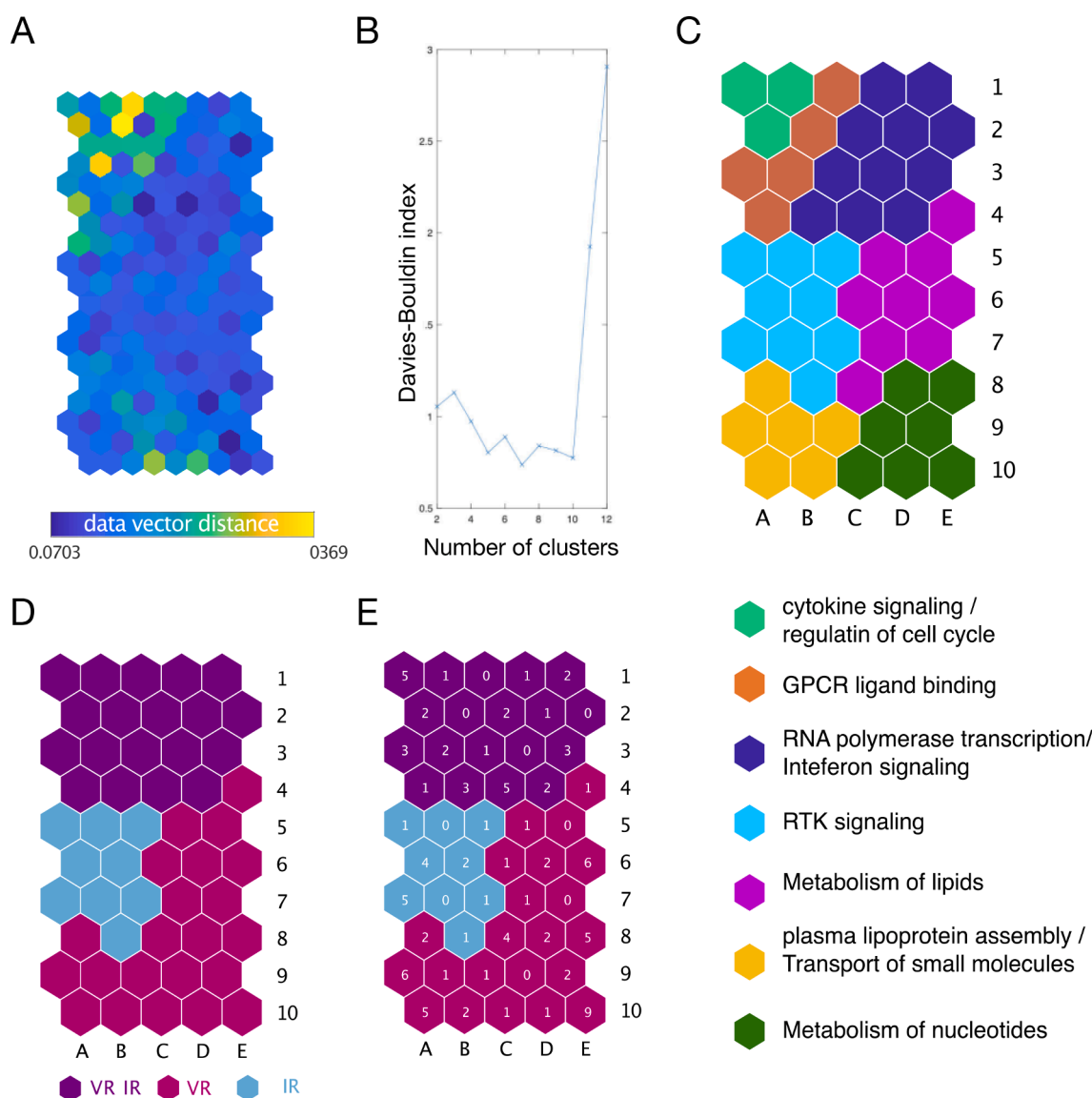
the MoA of drugs implicated for the treatment of SARS. We generated an F1 score matrix (for 102 drugs and 219 pathways in which drug target proteins are overrepresented, [Supplementary Table 3](#)) and then we investigated the MoA of identified drugs using SOM (Fig. 7).

We initially generated a 219 pathways U matrix (Fig. 7A) followed by separation of the 219 key pathways into seven clusters by a k-mean clustering algorithm with Davies-Bouldin (DB) index (Fig. 7B). The seven clusters were “Cytokine signaling or cell cycle”, “GPCR ligand binding”, “RNA polymerase transcription or Inteferon signaling”, “RTK signaling or platelet homeostasis”, “Metabolism of lipids”, “plasma lipoprotein assembly or Transport of small molecules”, “Metabolism of nucleotides or cell cycle” (Fig. 7C and [Supplementary Table 1](#)). Seven pathway clusters were mapped to 3 MoA categories: viral replication and immune responses, viral replication or immune responses, based on the Reactome pathway hierarchy and literature [36]. (Table 2, Fig. 7D). Finally, the SOM mapped the 102 drugs into each hexagon to specify the key pathways (the number of drugs per hexagon is shown in Fig. 7E). 15 drugs mapped to either the viral replication or the immune response

**Table 2**  
MoA of identified drugs for SARS.

Cluster	Pathway	MoA Category	Drugs
C1	cytokine signaling/regulatin of cell cycle	VR, IR	8
C2	GPCR ligand binding	IR	4
C3	RNA polymerase transcription/Inteferon signaling	VR, IR	22
C4	RTK signaling/platelet homeostasis	IR	15
C5	Metabolism of lipids	VR	16
C6	plasma lipoprotein assembly/Transport of small molecules	VR, IR	17
C7	Metabolism of nucleotides/cell cycle	VR	20

MoA pathways and notably, four out of the 21 drugs in the viral replication MoA category are in clinical trials for the treatment of COVID-19 [33], as are two of the drugs that mapped to the immune responses pathway (Fig. 7E, [Supplementary Table 2](#)).



**Fig. 7.** Predicted mechanism of actions of 102 candidate drugs for SARS. (A) U-matrix is shown of the trained unsupervised SOM used to analyze the relationship between the 102 drugs and the 219 key pathways. (B) 219 key pathways were separated into seven clusters by a k-mean clustering algorithm with Davies-Bouldin (DB) index. (C) SOM component maps (Fig. S4) are used to map the 219 pathways to the seven clusters, which are listed below the map along with their associated colour. (D) The seven clusters are mapped to three main therapeutic mechanisms of action using pathway analysis based on the Reactome pathway hierarchy information. (E) The SOM then assigns the 102 drugs to each relevant hexagon (neuron) to indicate a relevant MoA per drug.

#### 4. Discussion

In this study, we apply a network-based methodology for the systematic identification of approved drugs that can be potentially repurposed for the treatment of SARS and COVID-19. This builds on our previous study using the similar approach to identify drugs that might provide potential treatments for COVID-19 [7]. The integration of drug–target networks, coronavirus–host interactions, coronavirus-derived proteomics from human cell lines, and human protein–protein interaction networks has enabled us to establish essential disease-specific networks; and analysis of this network allowed us to identify key proteins associated with coronavirus infectious disease. Through the Artificial Neural Networks method, the MoA of the identified drugs was analysed to infer how the identified drug might impact SARS or COVID-19 biology.

Our analysis identifies 196 approved drugs for COVID-19 and 102 approved drugs for SARS, along with their MoA (Supplementary Table 2). To validate the predictive power of our approach, we identified the number of drugs predicted from our analysis that are currently in clinical trials for COVID-19. Of the 196 drugs identified for COVID-19, 44 were identified as being in clinical trials for COVID-19 patients, indicating that our methodology is statically significant. An important feature of our analyses is that only approved drugs were used. This allows for rapid advancement of the most promising of the 152 drugs which have not yet entered clinical trials for COVID-19. In the case of SARS, none of the identified drugs were found to be in clinical trials, so the statistical validation of the predicted drugs could not be carried out.

The list of 196 drugs identified as potential treatments for COVID-19 has a 40% overlap with the drugs identified in our previous COVID-19 study [7]. Although the DIPs and DEPs have only a 7% and 9% overlap, respectively, between the two studies (possibly owing to the proteomics datasets being from different cell types (kidney in the previous study [7] and lung in the current study)), the network-based methodology robustly identified drugs and biological pathways (more than 80% of similarity) in our two studies.

By constructing and analysing the CIP network, which includes an added hidden layer of proteins between DIP and DEP, we were able to find key target proteins for COVID-19 treatment. For example, among the predicted drugs we identified for SARS and COVID-19, the most frequently targeted protein is prostaglandin-endoperoxidase synthase 2 (PTGS2), a protein that is involved in inflammation and the generation of mitochondria [43]. PTGS2 is in the hidden layer of the CIP network and was selected as a key protein through centrality analysis and RWR. Of the 196 drugs identified for COVID-19, 49 drugs target PTGS2 and 13 of them are in clinical trials; this includes dexamethasone, which in the RECOVERY trial [44] reduced deaths from COVID-19 by about a third. Two other proteins in the hidden layer, AKT1 and nitric oxide synthase 3 (NOS3), which are necessary for viral synthesis [7,45,46] were also targeted by a number of the drugs identified for repurposing in this study. AKT1 is targeted by 19 drugs and NOS3 is targeted by 18 drugs. Therefore, through the construction and analysis of a CIP network consisting of DIPs, DEPs and a hidden layer of proteins, we were able to predict drugs that target key proteins associated with viral replication and immune responses during infection, some of which are already proven treatments for patients severely ill with COVID-19.

One limitation of our analyses is the use of undirected protein–protein network analysis, meaning it is unknown whether the predicted drug inhibits or accelerates coronavirus infection. For example, flucytosine was predicted to be associated with viral replication in this study and in our previously published study of viral replication, flucytosine accelerated rather than inhibited virus replication [7]. This limitation can be solved by using a disease-specific directional network, which could be generated using time-series transcription data, phosphorylation data and a human protein interaction network, such as a regulatory network and a protein–protein network. If the methodology we have described in this paper was used on a disease-specific directional

network, it could be used to predict drugs associated only with the suppression of virus infection.

#### 5. Conclusion

This study presents a robust and integrated network-based methodology to rapidly identify clinically approved drugs that can potentially be repurposed to treat SARS and COVID-19. The identified compounds have fully annotated biological mechanisms and therapeutically relevant information, which should accelerate the entrance of these drugs into clinical trials for severely ill patients with COVID-19. Furthermore, our unbiased, systematic, data-driven computational approach should be useful in identifying drugs that can be repurposed to treat future novel viral outbreaks.

#### CRediT authorship contribution statement

**Woochang Hwang:** Conceptualization, Methodology, Data curation, Software, Validation, Writing – original draft. **Namshik Han:** Conceptualization, Methodology, Writing – review & editing, Supervision.

#### Declaration of Competing Interest

N.H. is a cofounder of KURE.ai and an advisor at Biorelate, Promatix, and Standigm. The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Acknowledgements

The authors would like to thank Alison Schuldt and Nicola McCarthy for their helpful advice and comments and Méabh MacMahon for assistance editing the manuscript. NH and WH are funded by LifeArc.

#### Data statement

The authors confirm that the data supporting the findings of this study are available.

Software and data are available at: [https://github.com/wchwa/ng/Method\\_Pancorona](https://github.com/wchwa/ng/Method_Pancorona).

#### Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.jymeth.2021.11.002>.

#### References

- [1] WHO, Summary of probable SARS cases with onset of illness from 1 November 2002 to 31 July 2003, 2003.
- [2] WHO, Coronavirus disease (COVID-2019) situation report, 2021.
- [3] D.M. Morens, A.S. Fauci, Emerging pandemic diseases: how we got to COVID-19, *Cell* 182 (5) (2020) 1077–1092, <https://doi.org/10.1016/j.cell.2020.08.021>.
- [4] A. Stukalov, V. Girault, V. Grass, O. Karayel, V. Bergant, C. Urban, D.A. Haas, Y. Huang, L. Oubraham, A. Wang, M.S. Hamad, A. Piras, F.M. Hansen, M.C. Tanzer, I. Paron, L. Zinzula, T. Engleitner, M. Reinecke, T.M. Lavacca, R. Ehmann, R. Wölfel, J. Jores, B. Kuster, U. Protzer, R. Rad, J. Ziebuhr, V. Thiel, P. Scaturro, M. Mann, A. Pichlmair, Multilevel proteomics reveals host perturbations by SARS-CoV-2 and SARS-CoV, *Nature* 594 (7862) (2021) 246–252, <https://doi.org/10.1038/s41586-021-03493-4>.
- [5] D.E. Gordon, G.M. Jang, M. Bouhaddou, J. Xu, K. Obernier, K.M. White, M. J. O'Meara, V.V. Rezelj, J.Z. Guo, D.L. Swaney, T.A. Tummino, R. Hüttenhain, R. M. Kaake, A.L. Richards, B. Tutuncuoglu, H. Foussard, J. Batra, K. Haas, M. Modak, M. Kim, P. Haas, B.J. Polacco, H. Braberg, J.M. Fabius, M. Eckhardt, M. Soucheray, M.J. Bennett, M. Cakir, M.J. McGregor, Q. Li, B. Meyer, F. Roesch, T. Vallet, A. Mac Kain, L. Miorin, E. Moreno, Z.Z.C. Naing, Y. Zhou, S. Peng, Y. Shi, Z. Zhang, W. Shen, I.T. Kirby, J.E. Melnyk, J.S. Chorba, K. Lou, S.A. Dai, I. Barrio-Hernandez, D. Memon, C. Hernandez-Armenta, J. Lyu, C.J.P. Mathy, T. Perica, K.B. Pilla, S. J. Ganesan, D.J. Saltzberg, R. Rakesh, X.i. Liu, S.B. Rosenthal, L. Calviello, S. Venkataramanan, J. Liboy-Lugo, Y. Lin, X.-P. Huang, Y. Liu, S.A. Wankowicz,

